

NPS55-83-007

NAVAL POSTGRADUATE SCHOOL, Monterey, California



SIMTBED: A GRAPHICAL TEST BED FOR ANALYZING
AND REPORTING THE RESULTS OF A
STATISTICAL SIMULATION EXPERIMENT

by

P. A. W. Lewis
E. J. Orav
H. W. Drueg
D. G. Linnebur
L. Uribe

April 1983

Approved for public release; distribution unlimited

Prepared for:
Chief of Naval Research
Arlington, VA 22217

NAVAL POSTGRADUATE SCHOOL
Monterey, California

Rear Admiral J. J. Ekelund
Superintendent

David A. Schrad
Provost

This work was supported in part by the Office of Naval Research under
Grant NR-42-284.

Reproduction of all or part of this report is authorized.

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER NPS55-83-007	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) SIMTBED: A GRAPHICAL TEST BED FOR ANALYZING AND REPORTING THE RESULTS OF A STATISTICAL SIMULATION EXPERIMENT		5. TYPE OF REPORT & PERIOD COVERED Technical
		6. PERFORMING ORG. REPORT NUMBER
7. AUTHOR(s) P. A. W. Lewis D. G. Linnebur E. J. Orav L. Uribe H. W. Drueg		8. CONTRACT OR GRANT NUMBER(s)
9. PERFORMING ORGANIZATION NAME AND ADDRESS Naval Postgraduate School Monterey, CA 93940		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS 61153N: RR014-05-01 N0001483WR30026
11. CONTROLLING OFFICE NAME AND ADDRESS Chief of Naval Research Arlington, VA 22217		12. REPORT DATE April 1983
		13. NUMBER OF PAGES 42
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		15. SECURITY CLASS. (of this report) Unclassified
		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited.		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number)		
SIMTBED Statistical Simulation Graphics Serial Correlation Coefficient Simulation Gamma Distribution Regression-Adjusted Graphics		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number)		
A graphical test bed in which the results of a simulation experiment can be reported and analyzed is described. The test bed is based on the regression adjusted graphics and estimation methodology developed by Heidelberger and Lewis for regenerative simulation. From the graphics and associated numerics, the experimenter can summarize and see simultaneously relative properties, such as bias, normality and standard deviation, of several estimators of a characteristic of a population for up to 8 sample sizes. The evolution of these properties with sample size is also displayed. The graphics is		

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

supported on a line printer to make it and the program portable. The technique is illustrated by two examples, one concerning the effects of changes in data distribution on the behavior of the estimated lag one serial correlation coefficient and the other concerning the relative properties of several estimators of a Gamma distribution.

S/N 0102- LF- 014- 6601

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

SIMTBED: A graphical test bed for analyzing and reporting the results
of a statistical simulation experiment

P.A.W. Lewis E.J. Orav H.W. Drueg

D.G. Linnebur

L. Uribe

Naval Postgraduate School
Monterey, California
93940

U.S. Marine Corps
Washington, D.C.

Computer
Applications
Salinas, CA
93908

ABSTRACT

A graphical test bed in which the results of a simulation experiment can be reported and analyzed is described. The test bed is based on the regression adjusted graphics and estimation methodology developed by Heidelberger and Lewis for regenerative simulation. From the graphics and associated numerics, the experimenter can summarize and see simultaneously relative properties, such as bias, normality and standard deviation, of several estimators of a characteristic of a population for up to 8 sample sizes. The evolution of these properties with sample size is also displayed. The graphics is supported on a line printer to make it and the program portable. The technique is illustrated by two examples, one concerning the effects of changes in data distribution on the behavior of the estimated lag one serial correlation coefficient and the other concerning the relative properties of several estimators of a Gamma distribution.

1.0 Introduction

SIMTBED is a graphical display program that can be used via a simulation on a digital computer to (i) explore the distribution of a statistical estimator for a given sample size, (ii) to compare the properties of that distribution when the estimator is calculated for various sample sizes, and (iii) to contrast those properties under different estimation conditions. Those conditions are controlled by the experimenter but, most commonly, they will entail competing estimation procedures (e.g. maximum likelihood versus methods of moments, or jackknifed versus not jackknifed). The program is flexible enough to accommodate the imagination of most users and, in one of the examples, we also consider the effects of changes in the underlying distribution of the data.

One salient feature of the program is that it uses the same batch of simulated random variables (e.g. Normals) to explore the properties of all the estimators at various sample sizes. This is done for economy of computer time and could be important on slow computers; the price paid is that the analytical analysis provided by SIMTBED of its graphical output is performed on correlated samples.

To use the program it is necessary only to define the optional input parameters, supply the simulated random variables, and provide the Fortran functions which, when passed the data and subsample size, transform (if desired) the data subsample and compute the desired statistics. SIMTBED itself will subdivide and feed the data properly into the functions, produce boxplots and summary statistics, and compute regressions for the mean and variance of each estimator based on inverse subsample size. Up to three estimators can be used with the option to

produce equally scaled graphs for all the statistics.

The features of the program are more easily demonstrated by example rather than explanation and so we will proceed directly to two applications. The first application refers back to a simulation study done by Cox (1966) looking at the behavior of the estimated first order serial correlation coefficient, Fisher's z-transform of the estimated correlation, and the 2-fold jackknifed estimate of the correlation for i.i.d. Normal(0,1), $\chi^2(1)$ and Lognormal(0,1) data. The jackknife was originally proposed by Quenouille (1948) for the purpose of removing bias from the correlation estimate. The second application considers the problem of estimating the shape parameters for a highly skewed Gamma(.25) and a nearly Normal Gamma(5.0) sample using m.l.e., method of moments, 4-fold jackknifed m.l.e., and 4-fold jackknifed method of moments as the competing estimators.

Technical details concerning the SIMTBED software, not essential to interpreting and appreciating the output, can be found in Linnebur (1982), and an application to the analysis of output in a regenerative simulation can be found in Heidelberger and Lewis (1981).

2.0 Calculation of the First Serial Correlation Coefficient

It is known that for an independent sample from a population with finite variance, the distribution of the serial correlation coefficient (Anderson and Walker, 1964) is asymptotically Normal with mean zero and variances $1/n$, where n is the sample size. If the population is i.i.d Normal then the bias is exactly $-1/n$. Since those asymptotic properties are frequently used as approximations in tests of significance, it is important to know how valid the approximation would be in small samples

from a variety of distributions. We will look at that question in the next two sections and then go on to consider two alternative measures of correlation, Fisher's z-transform and the 2-fold jackknifed estimate of the correlation. Their ability to reduce bias and/or induce Normality will be examined against other changes in the distribution of the estimators, particularly variance inflation. A simulation study, without graphics, of some of these problems was conducted by Cox (1966). He did not consider the jackknifed estimate.

2.1 SIMTBED Output for Serial Correlation

Figure 1(a) shows the simulated distribution and sample properties of the serial correlation coefficient estimate

$$r_n = \frac{n \sum_{j=1}^{n-1} (X_j - \bar{X}_1) (X_{j+1} - \bar{X}_n)}{(n-1) \sum_{j=1}^n (X_j - \bar{X}_0)^2},$$

$$\text{where } \bar{X}_0 = \sum_{j=1}^n X_j/n, \quad \bar{X}_1 = \sum_{j=1}^{n-1} X_j/(n-1), \quad \text{and } \bar{X}_n = \sum_{j=2}^n X_j/(n-1)$$

for various sub-sample sizes $n=n_i$. This definition matches that used by Anderson and Walker (1964). We consider first subsamples of size $n_1 = 10$, and then of size $n_2 = 20$, $n_3 = 30$, $n_4 = 40$, $n_5 = 50$, $n_6 = 75$, $n_7 = 100$ and $n_8 = 150$, successively. For each subsample size the input sample of $N = 5000$ simulated $\text{Normal}(0,1)$ random variables is divided into as many full subsamples of size n_i as possible, and the serial correlation is computed for each of the $\lfloor N/n_i \rfloor$ subsamples of size

n_1 . The entire procedure is then replicated $M = 10$ times, each time with a new simulated sample of $N = 5000$ Normal(0,1) variables.

After all M replications have been run, all the estimates of serial correlation for each subsample size are grouped together and their simulated distribution is presented via a boxplot and summary statistics (see e.g. Fig. 1(a)). The boxplot follows the standards discussed in Mosteller and Tukey (1977) with the median denoted by a + within the box, the mean by a * within the box, the outliers by 0's, and the far outliers by *'s beyond the whiskers. The summary statistics include the sample mean, sample standard deviation, estimated standard deviation of the sample mean (i.e. sample standard deviation/sqrt($M \lfloor N/n_1 \rfloor$)), sample skewness and sample kurtosis of the correlation estimates.

Looking at the output, the first (leftmost) boxplot in the graph in Figure 1(a) shows the distribution of

$$(\# \text{ Replications}) \times \left\lfloor \frac{(\text{Simulation Sample Size})}{(\text{Subsample Size})} \right\rfloor = 10 \times \left\lfloor \frac{5000}{10} \right\rfloor = 10 \times 500 = 5000$$

estimates of serial correlation from independent subsamples of size $n_1 = 10$. Summary statistics for the boxplot can be found below the graph in the column labeled "Subsample Size 10", so that the average serial correlation is $-.1074$, and the estimated standard deviation is $.2996$. The estimated standard deviation of the serial correlation estimate is $.2996/\sqrt{(5000)} = .00424$. Recall that this refers to correlation estimates based on subsamples of size 10.

Since the X-axis of the graph represents subsample size, the last (rightmost) boxplot shows the distribution of

$$10 \times \left\lfloor \frac{5000}{150} \right\rfloor = 10 \times 33 = 330$$

estimates of serial correlation from independent subsamples of size $n_g = 150$. Although the 330 estimates are independent of each other, they are not independent of the 5000 estimates that comprise the first boxplot since the same data (divided and processed in different ways) was used for both. Summary statistics show that the average correlation has dropped to $-.007372$, indicating the fall off in bias, and the standard deviation has dropped to $.07822$, indicating the greater precision with which the correlation can be estimated when 150 points, rather than 10, are available.

In order to quantify the changes that are occurring in the mean and variance of the distribution of the estimator as subsample size changes, SIMTBED performs two types of regressions. The first regression is on the averages and is done after each replication, using the average serial correlation for that replication, \bar{r}_{n_i} , as the dependent variable. Inverse powers of the subsample size serve as the independent variables. For Figure 1(a) the degree of the regression was chosen to be $D=3$ so, for each replication, the equations we attempt to fit by least squares are:

$$\bar{r}_{n_i} = a_0 + \frac{a_1}{n_i} + \frac{a_2}{n_i^2} + \frac{a_3}{n_i^3} \quad \text{for } i = 1, 2, \dots, 8.$$

This form anticipates the general asymptotic expansion

$$E(\hat{\theta}(n)) = \theta + \frac{\alpha_1}{n} + \frac{\alpha_2}{n^2} + \frac{\alpha_3}{n^3} + \dots$$

which holds true in the current situation with $\theta = 0$ and (in the Normal case) $\alpha = -1$ (see Cramer, 1948, for general results of this type).

Values of a_0 , a_1 , a_2 , and a_3 are calculated after each replication, averaged across the M replications to get \bar{a}_0 , \bar{a}_1 , \bar{a}_2 , and \bar{a}_3 , and then the averages are reported below the summary statistics on the line "Mean of Regression on Averages - Coefficients". We find that

$\bar{a}_0 = -.000272$ and $\bar{a}_1 = -1.03074$, both close to their respective theoretical counterparts of zero and -1.

Because we have 10 replications and therefore 10 independent values of each of a_0 , a_1 , a_2 , and a_3 , we can also estimate the variances and standard deviations of \bar{a}_0 , \bar{a}_1 , \bar{a}_2 , and \bar{a}_3 across replications. These values are presented on the two lines immediately below the coefficients. For instance, the estimated s.d. of the estimate $\bar{a}_0 = -.000272$ of a_0 is .003892.

The regression line for the mean value of the estimator is presented visually in the graph as a dotted curve. The estimated asymptote (i.e. \bar{a}_0) is printed with a dashed line wherever it does not coincide with the regression line. Bias, therefore, can be viewed as the difference between those two lines.

The second regression referred to above is done after all replications have been run and the variances of the estimators at each subsample size have been calculated. (Note that the standard deviations, not the variances, are presented in the summary statistics.) It should be recalled from previous discussion that these variances, as well as all measures in the summary statistics, are based on the grouping together of the serial correlations from all replications, at each subsample size. This is in contrast to the the procedure for the regression on the means, where average correlations are computed for each subsample size for each replication. In the case of the variances, we have 8 equations:

$$\hat{\text{Var}}(r_{n_i}) = \frac{b_0}{n_i} + \frac{b_1}{n_i^{3/2}} + \frac{b_2}{n_i^2} + \frac{b_3}{n_i^{5/2}}, \quad i = 1, 2, \dots, 8,$$

which we fit by least squares in order to estimate the coefficients β_0 , β_1 , β_2 , and β_3 in the presumed asymptotic expansion

$$\text{Var } (\hat{0}(n)) = \frac{\beta_0}{n} + \frac{\beta_1}{n^{3/2}} + \frac{\beta_2}{n^2} + \frac{\beta_3}{n^{5/2}} + \dots$$

This expansion holds for the variance of the estimated serial correlation coefficient for independent data. Usually it will be β_0 in which we are most interested since β_0 is used in computing asymptotic relative efficiencies of estimators. For independent data with finite variance, we know that $\beta_0 = 1$. The computed values of b_0 , b_1 , b_2 , and b_3 , are presented on the line labeled 'Regression on Variance - Coefficients'. Notice that $b_0 = .7438$ is close to the theoretical value of 1.

The final two numbers on Figure 1(a), YMIN and YMAX, simply show the scale of the vertical axis. Because the SIMTBED program option to put Figures 1(a), 1(b) and 1(c) on the same scale was in effect, it may be that no boxplot in a given Figure (eg. Figure 1(b)) requires the full range of Y-values.

In order to produce Figure 1(b), the Normal(0,1) data that went into Figure 1(a) was squared to create longer tailed $\chi^2(1)$ random variables. The output is entirely analogous to that for Figure 1(a). Similarly, for Figure 1(c), the Normal(0,1) data was exponentiated in order to create Lognormal(0,1) data and to produce analogous graphical output. The indication is that the distribution of the sample serial correlation is robust with respect to the population distribution.

The features of the SIMTBED output will become clearer when they are associated with the various properties of the correlation estimator. First, however, a few technical comments concerning the regressions are necessary.

2.2 Some Comments on the Regressions

Two types of problems, numerical and statistical, can occur when

attempting to fit the two sets of regression equations presented in Section 2.1.

First, there is the question of numerical stability when the independent variables, $\{1, n_i^{-1}, n_i^{-2}, n_i^{-3}\}$ or $\{n_i^{-1}, n_i^{-3/2}, n_i^{-2}, n_i^{-5/2}\}$ decrease geometrically. If we attempt to form $X^T X$, where X is the respective design matrix and X^T is the transpose of X , we get values that range from 8 (assuming 8 subsample sizes) to $\sum_{i=1}^8 n_i^{-6}$ for the regression on the means, and $\sum_{i=1}^8 n_i^{-2}$ to $\sum_{i=1}^8 n_i^{-5}$ for the regression on the variances. Experience has shown that attempts to solve systems with such extremes in the $X^T X$ matrix produce erroneous results. Consequently, SIMTBED scales the design matrices by multiplying each entry of X by $\text{Max}(n_i)$ raised to the proper power so that no entry becomes too small. The standard Choleski decomposition (see Dahlquist, Bjorack and Anderson, 1974) is then used to fit the equations, and the coefficients are properly rescaled before they are reported. This procedure produces numerically reliable results.

The second problem concerns the breakdown of statistical assumptions in our regression models. It has already been pointed out in Section 2.1 that the two sets of dependent variables:

(1) the $\bar{\theta}(n_i)$ when considering the regression on the means;

$$(2) \text{ the } s^2(n_i) = \frac{M \left[\sum_{j=1}^{N/n_i} (\hat{\theta}_j(n_i) - \bar{\theta}(n_i))^2 \right]}{M \left[N/n_i \right]},$$

where $\bar{\theta}(n_i)$ is the mean across the M replications, when considering the regression on the variances,

Table 1

Entries in the table are the estimated correlations between the estimated variances of the r_{n_i} at different subsample sizes: $\text{Corr}(s^2(r_{n_i}), s^2(r_{n_j}))$ for $i=1, \dots, 8, j=1, \dots, 8$.

i \ j	1	2	3	4	5	6	7	8
1	1.00	.49	.46	-.26	.18	-.17	.14	.01
2	.49	1.00	.40	.55	.11	.38	.38	-.03
3	.46	.40	1.00	.23	.23	.44	.21	.29
4	-.26	.55	.23	1.00	.42	.86	.57	.35
5	.18	.11	.23	.42	1.00	.71	.43	.59
6	-.17	.38	.44	.86	.71	1.00	.45	.53
7	.14	.38	.21	.57	.43	.45	1.00	.72
8	.01	-.03	.29	.35	.59	.53	.72	1.00

Recall that r_n is the estimated serial correlation for a simulated Normal(0,1) subsample of size n . Also, the estimated correlations shown above were computed using 10 values (replications) of $s^2(r_{n_i})$ and $s^2(r_{n_j})$ for each i and j .

Table 2

A comparison of the estimated variance of $s^2(r_{n_i})$ with the approximate theoretical variance of $s^2(r_{n_i})$ and with the approximately equivariant scaled versions, $n_i^{-.5} s^2(r_{n_i})$. All entries have been multiplied by 10^5 .

$n_i =$	10	20	30	40	50	75	100	150
$\hat{\text{Var}}(s^2(r_{n_i}))$.177	.150	.204	.079	.047	.031	.049	.022
Approx. Theoretical $\text{Var}(s^2(r_{n_i}))$.400	.200	.133	.100	.080	.053	.040	.027
$\hat{\text{Var}}(n_i^{-.5} s^2(r_{n_i}))$	1.77	2.99	6.12	3.18	2.33	2.33	4.88	3.35

The estimated variances of $s^2(r_{n_i})$ and $\sqrt{n_i} s^2(r_{n_i})$ were calculated using 10 independent replications of $s^2(r_{n_i})$.

are not independent over i since all are based on the same simulated data. The extent of the dependence is demonstrated by the correlation matrix in Table 1. Entries in that table show the estimated correlation between $s^2(n_i)$ and $s^2(n_j)$ for all i and j , where the estimation was done by repeating the SIMTBED experiment with 10 different batches of 50,000 simulated random variables. Since only 10 values went into each correlation calculation, the table is only accurate to within approximately $\pm 2/\sqrt{10} = .632$. We see some indication of positive correlation, especially when i and j are close, but the lack of independence is not severe enough to hurt the regression results for either the estimated means or variances significantly.

A second assumption, implicit in any regression, is that the dependent variables have equal variances. This condition holds true for the means, which can be shown to satisfy

$$\text{Var}(\bar{\theta}(n_i)) = \frac{M}{N}$$

independently of i . The estimated variances, however, are not equivariant and, if we assume the $\hat{\theta}_j(n_i)$ to be approximately Normally distributed so that

$$M \sum_{j=1}^{\lfloor N/n_i \rfloor} (\hat{\theta}_j(n_i) - \bar{\theta}(n_i))^2$$

is approximately proportional to a Chi-squared random variable, with $M \lfloor N/n_i \rfloor - 1$ degrees of freedom, we can compute

$$\text{Var}(s^2(n_i)) \cong \frac{2}{MNn_i - n_i^2}.$$

To correct this problem of unequal variances SIMTBED scales the $s^2(n_i)$ by $\sqrt{n_i}$ so that

$$\text{Var} (\sqrt{n_i} s^2(n_i)) \cong \frac{2}{MN - n_i} \cong \frac{2}{MN}$$

since $MN \gg n_i$. The design matrix is scaled accordingly and the values b_0 , b_1 , b_2 , and b_3 discussed in Section 2.1 are reported.

Table 2 shows the effects of the rescaling by presenting first the estimated variances of the $s^2(n_i)$, where the estimation is done by repeating SIMTBED for 10 batches of 50,000 simulated data points. These estimated variances decrease as n_i increases, closely paralleling the second line of Table 2 which has the approximate theoretical values (i.e. $2/(MNn_i - n_i^2)$). The final line of Table 2 shows the estimated variances of the rescaled $s^2(n_i)$, i.e. the $\sqrt{n_i} s^2(n_i)$, which, as expected and hoped, show a more constant variance with i .

Although future versions of SIMTBED will include more sophisticated regression routines and the ability to generate independent samples at each subsample size, the current version is quick, usable, and accurate for most situations.

2.3 Interpreting the Serial Correlation Results

Returning to Figure 1(a) which shows the simulated distribution of the serial correlation coefficient from independent, Normal(0,1) data, the following comments summarize the most striking features:

(a) The boxplots appear very symmetric at all subsample sizes with nearly equal numbers of outliers at either tail and with mean and median coincidental. This observation is confirmed by the estimates of skewness in the summary statistics. Kurtosis is mildly negative at small subsample sizes but, overall, asymptotic Normality seems to take hold rather quickly. Note however, that at $n_i = 10$ there are only 3 outliers in a

sample of size 5000. This is consistent with the estimated kurtosis of -0.424 , showing that the distribution is quite nonnormal.

(b) The average serial correlation is negative for small subsamples. This is demonstrated by the dotted regression curve which starts at approximately $-.10$ and levels off near 0 for subsamples greater than about 85. The dashed asymptote of $-.000272$ is very close to the theoretical value of 0 , and the mean values in the summary table closely reflect the bias of $-1/n$.

(c) The standard deviations of the simulated distributions are very close to the asymptotic values of $n_i^{-0.5}$, although the lead coefficient in the regression on the variances, $b_0 = .743756$, is not as close to the theoretical value of 1 as we would hope. When SIMTBED is repeated 10 times with 10 different batches of simulated data, we find an average value for b_0 to be 1.0604 , with a standard deviation for b_0 of $.307$. The estimation procedure for b_0 , therefore, remains valid, but the estimate itself is highly variable.

The agreement between the simulated and the theoretical, asymptotic values of the bias and variances was discovered previously by Cox (1966). SIMTBED has now allowed us to automatically look at a broader range of subsample sizes and to see, through boxplots and estimates of skewness and kurtosis, a fuller picture of any changes in the distribution of the estimator. We can be satisfied that estimates of serial correlation do behave approximately as $\text{Normal}(-1/n, 1/n)$ random variables when the underlying data is $\text{Normal}(0,1)$.

If the lead terms in the expansions of the mean and variance of the estimated correlation coefficient (ie. a_0 , a_1 , and b_0) had been unknown, we would also have a fairly good idea now of what they were.

When the underlying data is χ_1^2 , Figure 1(b) confirms Cox's observation that the bias is relatively unaffected but, for small subsamples, the standard deviation is smaller than the expected $n^{-0.5}$. Unlike Figure 1(a), there is a pronounced skewness in the boxplots in Figure 1(b) with many more outliers at the positive end, and with the mean higher than the median at the first four subsample sizes. The problem of suppressed variance seems cured at $n_7 = 100$ and $n_8 = 150$, but the skewness remains and could cause problems in tests of significance.

Figure 1(c), which is based on an underlying batch of simulated Lognormal(0,1) data, shows a slight exaggeration of the effects in Figure 1(b). The standard deviation is more suppressed and does not attain the theoretical level by $n_8 = 150$. The positive skewness is more pronounced and kurtosis does not approach the theoretical value of 0.

Overall, the effects of long-tailed data on the distribution of the serial correlation coefficient can be summarized as follows:

- (i) Bias is not significantly affected and remains at approximately $-1/n$.
- (ii) The variance of the distribution of the serial correlation coefficient is reduced by longer-tailed data.
- (iii) Positive skewness is created in the distribution.
- (iv) Kurtosis may become positive at large subsample sizes.
- (v) For long-tailed data (i.e. Lognormal), a subsample size of 150 is not large enough to insure asymptotic Normality.

2.4 SIMTBED Output for the z-Transform of the Correlation

Fisher's z-transform of the estimated correlation coefficient is defined by:

$$z_n = \frac{1}{2} \log \frac{1 + r_n}{1 - r_n} ,$$

where r_n is the estimated serial correlation presented in Section 2.1. The transformation is intended to make the distribution of the z_n more Normal than that of the r_n . When the same SIMTBED experiment described in Section 2.1 is run using z_n as the estimator instead of r_n , we get the results shown in Figures 2(a), 2(b) and 2(c). It should be noted that the scale of the boxplots here has been forced to be approximately comparable to the scale for the boxplots in Section 2.1. This is done by suppressing outliers that are more than 1.5 interquartile distances beyond the quartiles of the boxplot. If we had allowed the data to scale the boxplots, we would have seen a much wider range on the vertical axis because the z_n are not restricted to the limits of -1 to +1 and because there is one far outlier at -3.8. In this type of "reduced graphics", we still see the number of outliers that fall beyond the allowable range through the numbers at the ends of the boxplots, but we do not see their actual locations.

Figure 2(a) shows the distribution of the z-transformed correlation coefficients when we use simulated Normal(0,1) data. At each subsample size, the mean and standard deviation are close to the theoretical n^{-1} and $n^{-1/2}$ respectively. The skewness and kurtosis at subsample size $n_1 = 10$ are far from the theoretical Normal distribution values of 0 and 0, reflecting partly the one far outlier at -3.8 and partly the negative skew in the remainder of the z_{n_1} 's. For other subsample sizes, there is no strong evidence to contradict the assumption of approximate Normality.

The relationship between Figure 2(b) and 2(a) is similar to that between 1(b) and 1(a). Figure 2(b), which is based on simulated x_1^2

data, shows (a) bias that is the same as for the transformed correlations based on Normal data, (b) slightly suppressed variances, particularly at small subsample sizes and (c) positive skewness which persists at large subsample sizes. In addition, there are signs of positive kurtosis at small subsample sizes.

Figure 2(c) is based on Lognormal(0,1) data and shows high values of skewness and kurtosis at almost all subsample sizes. Approximate Normality seems an unwarranted assumption. In fact, the kurtosis is converging very slowly to its asymptotic value of 0.

In general, using the z-transform does not help with Normality assumptions, especially when dealing with long-tailed distributions.

2.5 SIMTBED Output for the 2-Fold Jackknife of the Correlation

The final Figures, 3(a), 3(b) and 3(c), deal with the 2-fold jackknife estimate of correlation. Again, the figures are reduced graphics with scaling comparable to that of the boxplots of Sections 2.3 and 2.4. To define the estimator, we start with a given subsample of size n , compute the serial correlation for the first $\lfloor n/2 \rfloor$ points and call it $r_1(n/2)$, compute the serial correlation for the second $\lfloor n/2 \rfloor$ points and call it $r_2(n/2)$ and compute the serial correlation for the entire subsample of n points and call it $r_0(n)$. Each computation follows the formula in Section 2.1. The three estimators are then combined to form two pseudo-values,

$$r_1^*(n) = 2r_0(n) - r_1(n/2)$$

and

$$r_2^*(n) = 2r_0(n) - r_2(n/2) ,$$

and the final jackknife estimator for that subsample is defined as

$$\tilde{r}(n) = \frac{r_1^{*(n)} + r_2^{*(n)}}{2} .$$

Although a jackknife estimator may have many favorable properties, we are concerned here primarily with its ability to remove bias, hopefully without inflating the variances of the estimator and/or inducing nonnormality.

Figure 3(a), based on simulated Normal(0,1) data, shows nearly complete removal of bias, even at small subsample sizes. The cost of the bias reduction is reflected in an increase of nearly 50% in the standard deviation of the correlation estimate for subsample size 10, and lesser relative increases at larger subsample sizes. There is also an indication of a positive skew for small subsample sizes, and the problem that the jackknife estimator need not fall into the -1 to +1 range which is desirable for a correlation coefficient estimate.

When using simulated χ_1^2 data as in Figure 3(b), or simulated Lognormal(0,1) data as in Figure 3(c), there is again no problem with bias. Variance inflation, though it exists at small subsample sizes, is not as large as when Normal(0,1) data is used. The distributions of the jackknifed correlations show very pronounced positive skews, however, as well as positive kurtosis. These two problems are worse for the longer-tailed Lognormal data.

Overall, the jackknife estimator is very successful at removing bias but the costs include variance inflation, which can be severe at small subsample sizes, plus increased positive skewness and kurtosis when the estimates are based on data from longer-tailed distributions.

2.6 Comparison of the Three Estimates of Correlation

For Normal(0,1) data, the distribution of the usual correlation coefficient displayed in Figure 1(a) behaves very much as theoretical

asymptotic calculations would predict, even at small subsample sizes. This makes it possible to correct for bias in the estimator and to perform tests of significance. Use of Fisher's z-transform, as illustrated in Figure 2(a) does not seem necessary since it does not significantly improve the approximate Normality of the estimator. The jackknife estimator in Figure 3(a) may be valuable if a direct, unbiased estimator is needed but the inflated variance of the jackknife estimator may limit the usefulness of the estimate as well as make any tests of significance too conservative.

When the underlying data comes from a longer-tailed distribution, the usual correlation coefficient in Figures 1(b) and 1(c) retains a predictable bias term, although the variance of its distribution is slightly depressed and the skewness and kurtosis becomes positive, even for subsamples as large as 150. This means that it is still possible to estimate the correlation accurately, but tests of significance fall on shaky assumptions of Normality. The z-transform in Figures 2(b) and 2(c) does little to firm up those assumptions and, in some cases, makes the situation worse. As in the case of Normal data, the 2-fold jackknifed correlation in Figures 3(b) and 3(c) is bias-free but follows a fairly nonnormal distribution which would invalidate significance testing.

All of the preceding observations and conclusions flow immediately from the nine Figures presented so far. Further studies could easily be done through SIMTBED, looking at larger subsample sizes, correlated data, i.e. $\rho \neq 0$ and alternative marginal distributions. For demonstration purposes, though, it is better to proceed to our second application.

3.0 Estimating the Shape Parameter for a Gamma Distribution

As a second application of SIMTBED, we will consider a problem

which has received much less statistical attention; asymptotic results are summarized in Cox and Lewis (1966, Ch.3) and Johnson and Kotz (1970, Ch.17). We want to estimate the shape parameter, K , for a Gamma distribution, where the Gamma density is given by

$$f(x) = \left(\frac{K}{\mu}\right)^K \frac{x^{K-1}}{\Gamma(K)} e^{-Kx/\mu} \quad x \geq 0; \quad K > 0; \quad \mu > 0$$

$$= 0 \quad x < 0.$$

Notice that the mean of this distribution is μ , not K/μ as in some differently parameterized versions of the Gamma density. For the data that will be simulated for use in SIMTBED we will use $K = 5$ and $\mu = 1$ and $K = 0.25$ and $\mu = 1$. The closer the mean of our estimate is to 5 or 0.25, the better (in terms of bias) is our estimation procedure. Other factors such as the variance and Normality of the estimator will of course also have influence in the determination of a preferred estimator; the bias and variance could be combined into m.s.e.

Section 3.1 will compare the commonly used maximum likelihood estimator which is mildly difficult to compute to the competing method of moments estimator which is very simple to compute. Both procedures result in asymptotically Normal estimators (Cramer, 1948) but the m.l.e. is usually preferred because of its favorable asymptotic relative efficiency (Cox and Lewis 1966). Through SIMTBED, though, we will see that for small subsamples the estimated variances of the two estimators of K are not as far apart as asymptotic results lead us to believe. In addition, the bias that appears in both estimators is smaller for the moment estimator.

In Section 3.2 we will use a four-fold jackknife of both the m.l.e. and moments estimators to successfully remove the bias. What is remarkable is that unlike the jackknifing of the serial correlation, there is

little or no cost in terms of variance inflation and nonnormality for the jackknifed moment estimator. When $K = .25$, we will see in Section 3.3 that the jackknifed m.l.e. dominates the other three estimators at all subsample sizes when considering the mean, variance, and Normality of the estimator.

3.1 Maximum Likelihood and Moment Estimators of K

Figure 4(a) is very similar in format to the figures that have already been presented for the correlation example except that:

(1) The estimator whose distribution is being displayed is the maximum likelihood estimator of K , the shape parameter of a Gamma(5) population. We denote the estimator, computed from a simulated subsample of size n , by $\hat{K}(n)$ and define it to be the solution of the equation:

$$n [\log \hat{K}(n) - \Psi(\hat{K}(n))] = n \log \sum_{i=1}^n X_i/n - \sum_{i=1}^n \log X_i ,$$

where the X_i are the simulated Gamma(5) random variables and $\Psi(\cdot)$ is the digamma function (Cox and Lewis, 1966).

(2) The eight subsample sizes which we will be looking at are $n_1 = 33$, $n_2 = 50$, $n_3 = 71$, $n_4 = 100$, $n_5 = 125$, $n_6 = 166$, $n_7 = 250$ and $n_8 = 500$. Note the difference between the n_i 's in the previous example and these n_i 's. Since these are larger we will not see much small sample detail, but we will see some of the asymptotic ($n = 500$) effects coming in.

(3) At each subsample size we will work with $M^* = 20$ independent replications of $N^* = 2500$ simulated Gamma(5) random variables, instead of the $M = 10$ replications of $N = 5000$ variables used previously.

The total number of independent simulated random variables across replications remains constant at the program maximum of 50,000. Hence, the boxplot at subsample size 50 in Figure 4(a) represents the distribution of $M^* \lfloor N^*/50 \rfloor = 1000$ estimates of $\hat{K}(50)$ just as the boxplot at subsample size 50 in Figure 1(a) represents the distribution of $M \lfloor N/50 \rfloor = 1000$ estimates of $r(50)$. As long as the product, $M \times N$, remains constant, the only effect of changing the number of replications is, up to rounding in N/n_i , to change the results in the regression on the averages. By using $M^* = 20$ and $N^* = 2500$, SIMTBED reports regression coefficients averaged over 20 replications, but, within each replication, the dependent variables are averages over just $\lfloor 2500/n_i \rfloor$ values of the estimator.

(4) The boxplots are presented using the reduced graphics option. In this option any extreme outliers (i.e. those beyond 1.5 interquartile distances) are included as a count at the tail of each boxplot. This option was chosen in order to give more graphical weight to the body of the distributions and the fall-off in the bias. Limited printer resolution makes it impossible to show details in the body and the tails of the distributions if there are many straggling outliers. In the case of very extreme outliers, no detail would be seen in the body of the boxplot without the reduced graphics option.

Figure 4(b) looks at the distribution of the moment estimator of K , the shape parameter of a $\text{Gamma}(K)$ population:

$$\tilde{K}(n) = (n-1) \bar{X}^2 / \sum_{i=1}^n (X_i - \bar{X})^2 ,$$

where $\bar{X} = \sum_{i=1}^n X_i / n$, n is the subsample size, and the X_i are the simulated $\text{Gamma}(5)$ random variables. The SIMTBED options and para-

meters mentioned in (2), (3) and (4) preceding are also in effect here.

The two figures, 4(a) and 4(b), show a very pronounced bias in both estimation procedures, although the moment estimator is slightly closer to the unbiased value of 5. As expected, the standard deviation of the m.l.e. is lower than that of the moment estimator although the relative difference at small subsample sizes, for instance 1.448 versus 1.482 at $n_1 = 33$, may not outweigh the increase in bias with the m.l.e. At larger subsample sizes, the relative difference is close to the theoretical asymptotic relative efficiency of .78 (i.e. .91 at $n_7 = 250$)

Both estimators also show distributions with positive skewness and kurtosis that decrease to the asymptotic 0 levels as subsample size increases. The asymptotics appear to take hold more quickly for the moment estimator than for the m.l.e.

In summary, SIMTBED shows that the m.l.e. is indeed better than the moment estimator in terms of variance, but not as good for small sample sizes as asymptotic results would lead us to believe. In the other areas of bias and asymptotic Normality, the moment estimator would have to be preferred.

3.2 4-Fold Jackknifed Estimators of K

Figures 4(c) and 4(d) show the distributions of the 4-fold jackknifed m.l.e. of K and 4-fold jackknifed moment estimator of K, respectively. A 4-fold jackknife estimator is similar to the 2-fold jackknife estimator described in Section 2.5 except that there are 4 pseudo-values that come out of dividing each subsample into fourths. More details can be found in Mosteller and Tukey (1977).

The purpose of the jackknife is to remove the conspicuous bias

observed in Figures 4(a) and 4(b). This goal is seen to be accomplished in Figure 4(c) and 4(d) and we can also note smaller values of skewness and kurtosis, indicating a quicker approach to asymptotic Normality. The skewness and kurtosis of the jackknifed moment estimator are the lowest, at small subsample sizes, among all estimators. The variance of the jackknifed moment estimator is also only slightly inflated, as is the variance of the jackknifed m.l.e..

All told, the jackknifed moment estimator, because of its lack of bias, small variance, and low skewness and kurtosis, would be the method of choice if estimation of K or significance testing was the goal.

3.3 Results for $K = 0.25$

In Figures 5(a), 5(b), 5(c) and 5(d) we show similar results to those discussed above for the case $K = 5.0$, but using $K = 0.25$. The fact (Cox and Lewis, 1966, Ch.3) that the m.l.e. estimate is much more efficient than the moment estimate is graphically illustrated. What is new is the effect of jackknifing: bias is reduced without the sacrifice of variance inflation or nonnormality.

Further comparisons and interpretations are similar to those done for the case $K = 5.0$, and are left to the reader.

4.0 Conclusions

Simply by providing SIMTBED with the desired estimators, we have been able (a) to explore in depth the effects of changes in data distribution and of different estimation procedures on the calculation of the serial correlation coefficient, and (b) to compare four different ways to estimate the shape parameter in a highly skewed Gamma population.

The graphics and numerical output combine to let us see and quantify distributional changes that occur as subsample size grows. We can see bias fall away, variance shrink, and skewness disappear as the estimator approaches asymptotic Normality. Terms in the asymptotic expansion of the mean and variance of the estimator are automatically calculated and can be used to compare different estimators.

Ease of use and portability, however, remain as SIMTBED's most important features, and will hopefully inspire users to try more diverse and extensive simulation experiments.

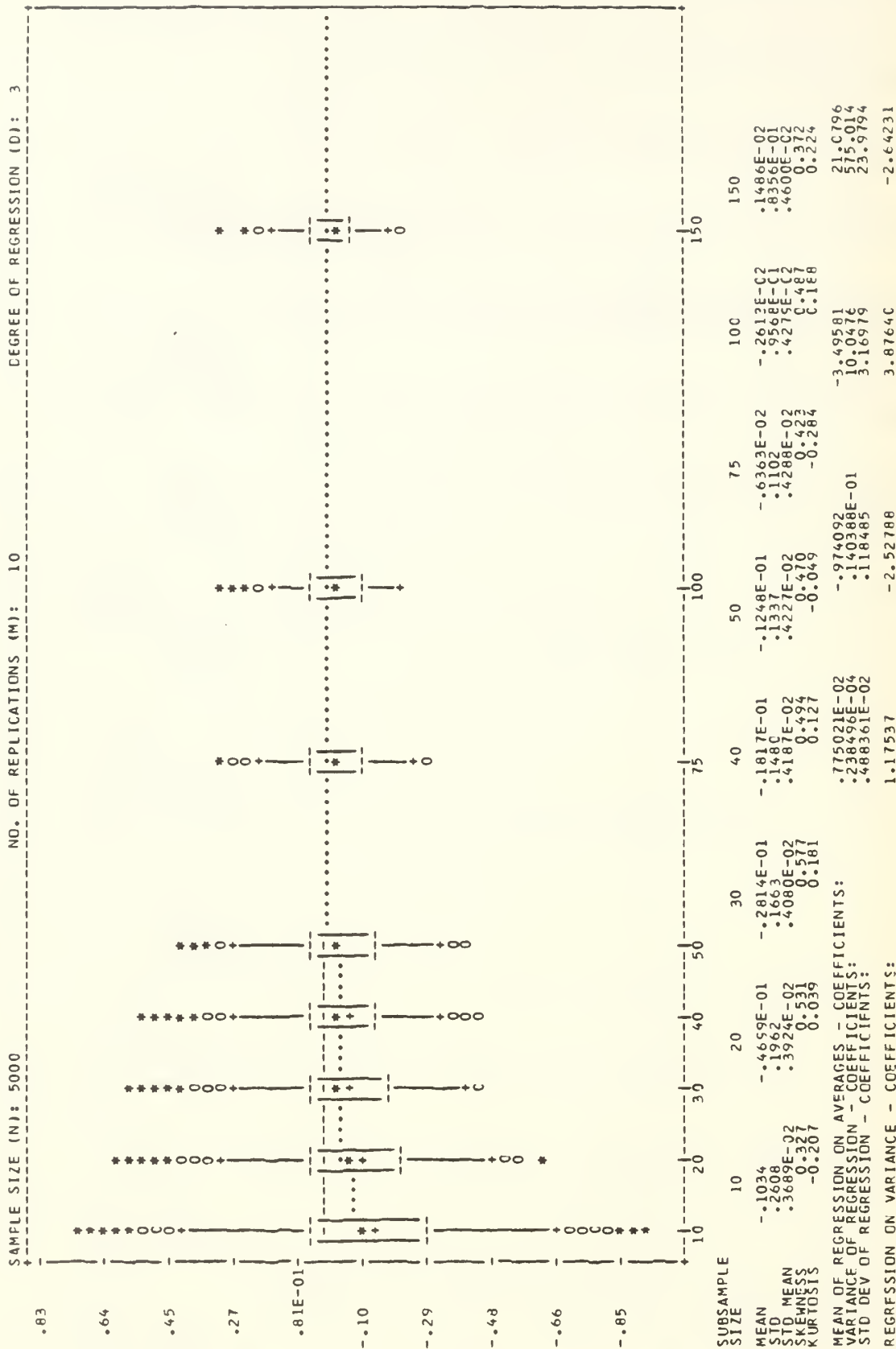
Other graphical displays besides running boxplots can be used; some alternatives are given in Lewis (1972), Heidelberger and Lewis (1981) and Devlin, Gnanadesikan and Kettenring (1981).

5.0 Availability

SIMTBED is at present only available in a version run on the IBM 3033. However, since the program uses standard FORTRAN and is independent of any software packages, conversion should be simple. Versions for VAX machines will be tested shortly.

References

- Anderson, T.W. and Walker, A.M., (1964). "On the Asymptotic Distribution of the Autocorrelation of a Sample from a Linear Stochastic Process", Ann. Math. Statist., 35, pp. 1296-1303.
- Cox, D.R., (1966). "The Null Distribution of the First Serial Correlation Coefficient", Biometrika, 53, pp. 523-626.
- Cox, D.R. and Lewis, P.A.W., (1966). The Statistical Analysis of Series of Events, Chapman and Hall, London, Ch. 6.
- Cramer, H. (1946). Mathematical Methods of Statistics, Princeton University Press, Princeton, New Jersey.
- Heidelberger, P. and Lewis, P.A.W. (1981). "Regression-Adjusted Estimates for Regenerative Simulations, with Graphics", Comm. of the A.C.M., 24, pp. 260-273.
- Johnson, N.L. and Kotz, S., (1970). Continuous Univariate Distributions - 1, Houghton Mifflin Company, Boston, Ch. 17.
- Linnebur, D.G., (1982). "A Graphical Test Bed for Analyzing and Reporting the Results of a Simulation Experiment", Master's Thesis, Naval Postgraduate School, Monterey, California.
- Quenouille, M.H., (1948). "Some Results in the Testing of Serial Correlation Coefficients", Biometrika, 35, pp. 261-267.



ESTIMATOR: ESTIMATES OF THE LAG ONE SERIAL CORRELATION COEFFICIENT FOR A CHI-SQUARE(1) SAMPLE

Figure 1(b)

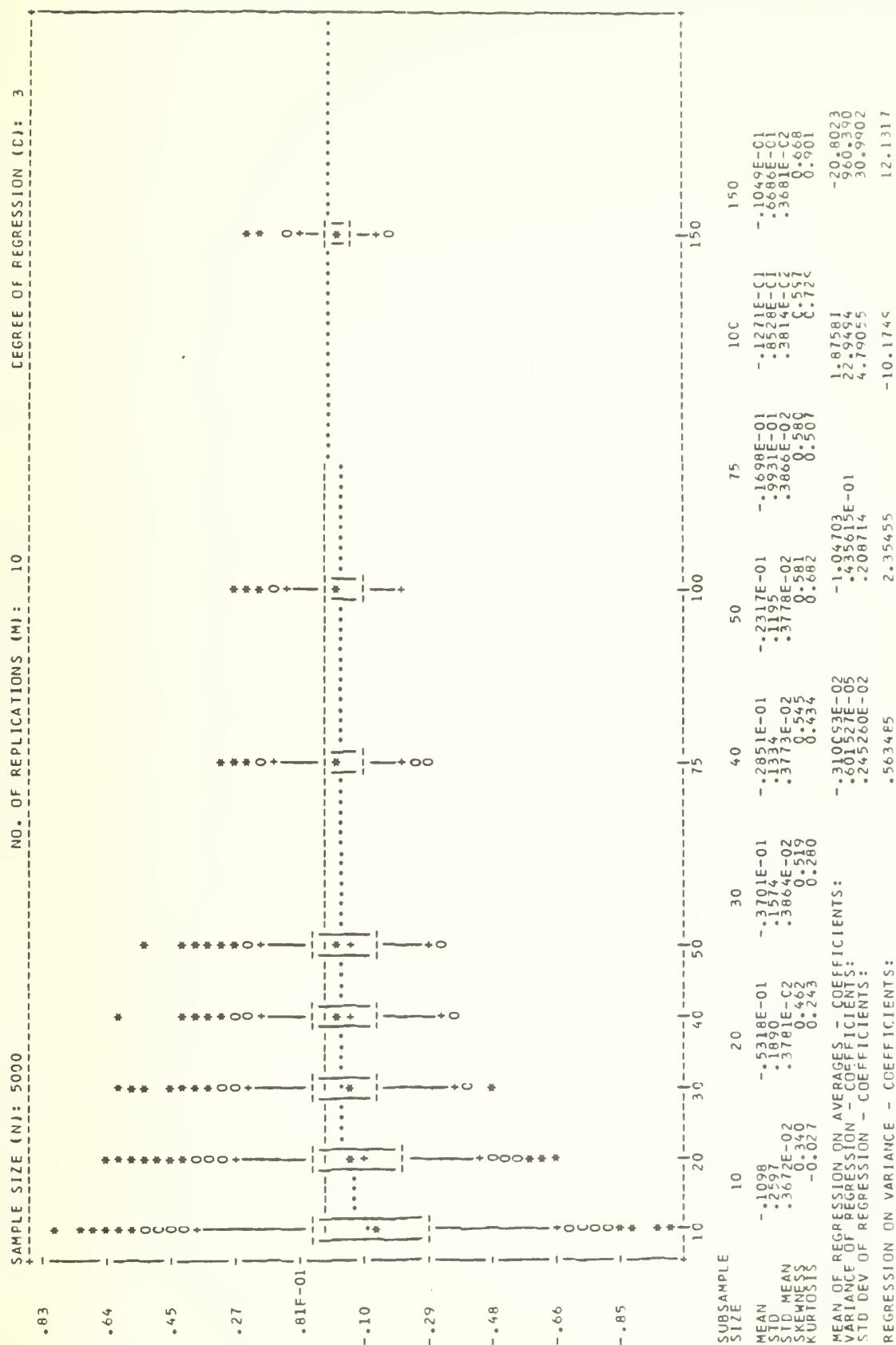


Figure 1(c)

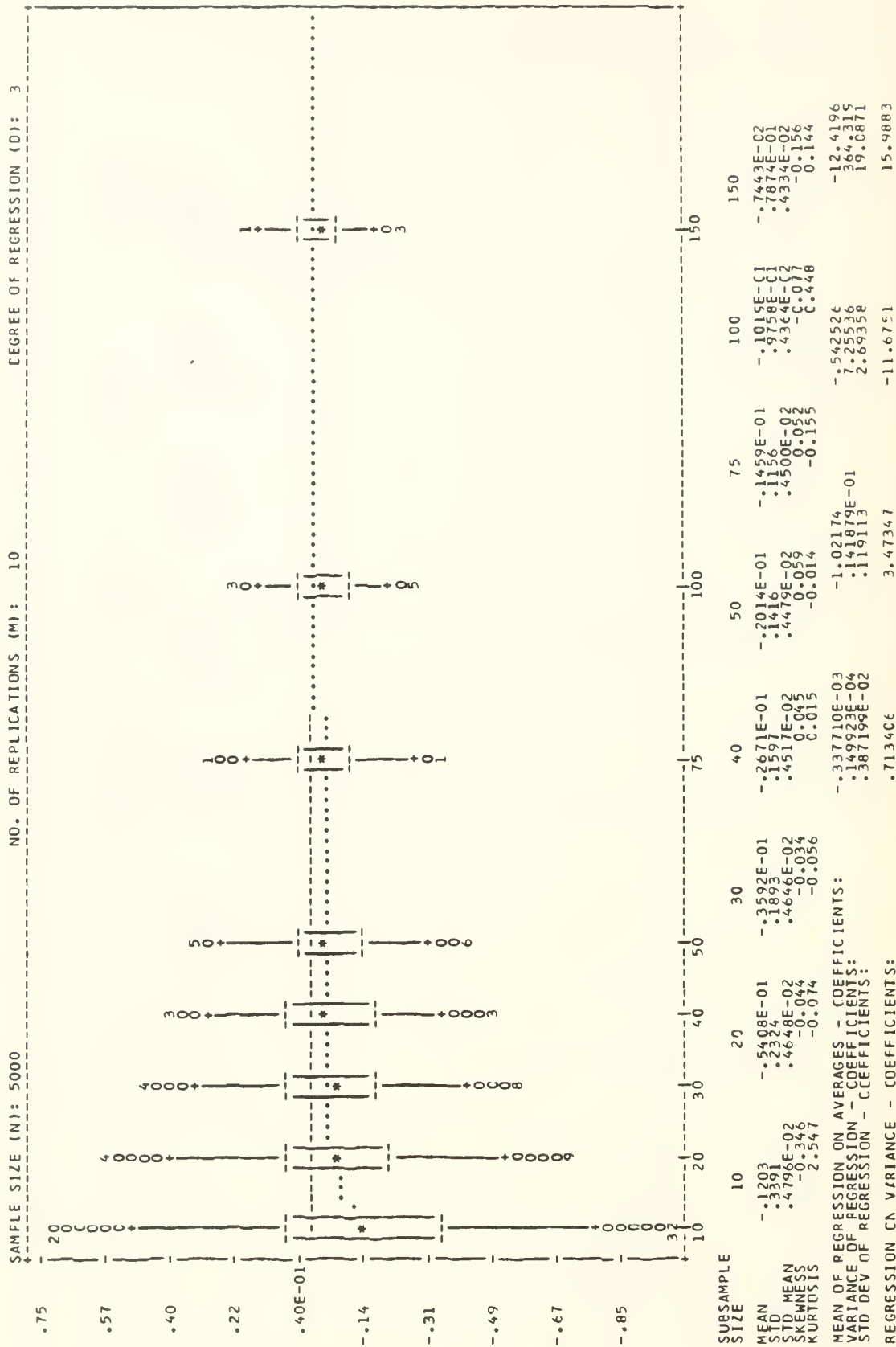
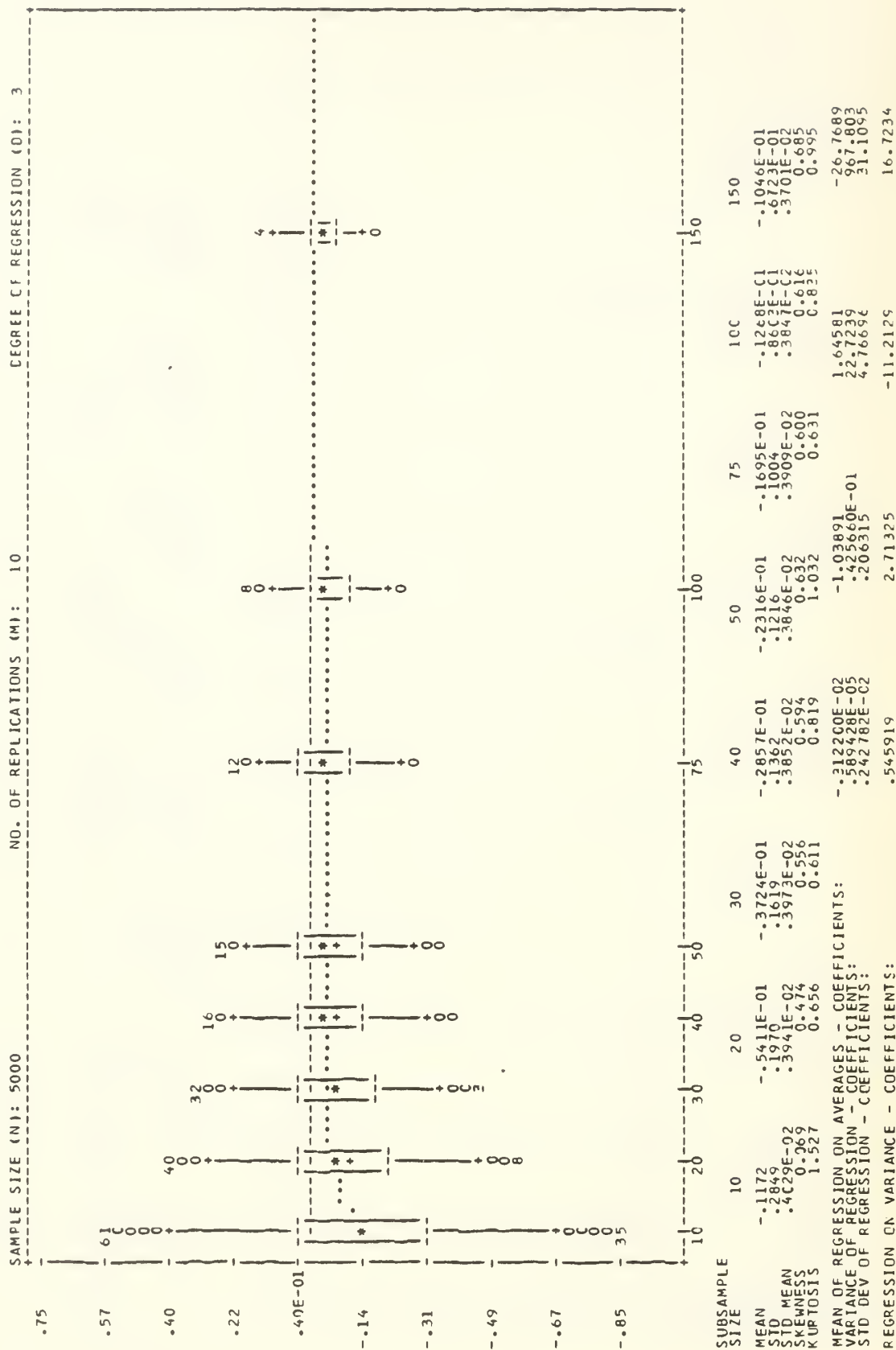


Figure 2(a)



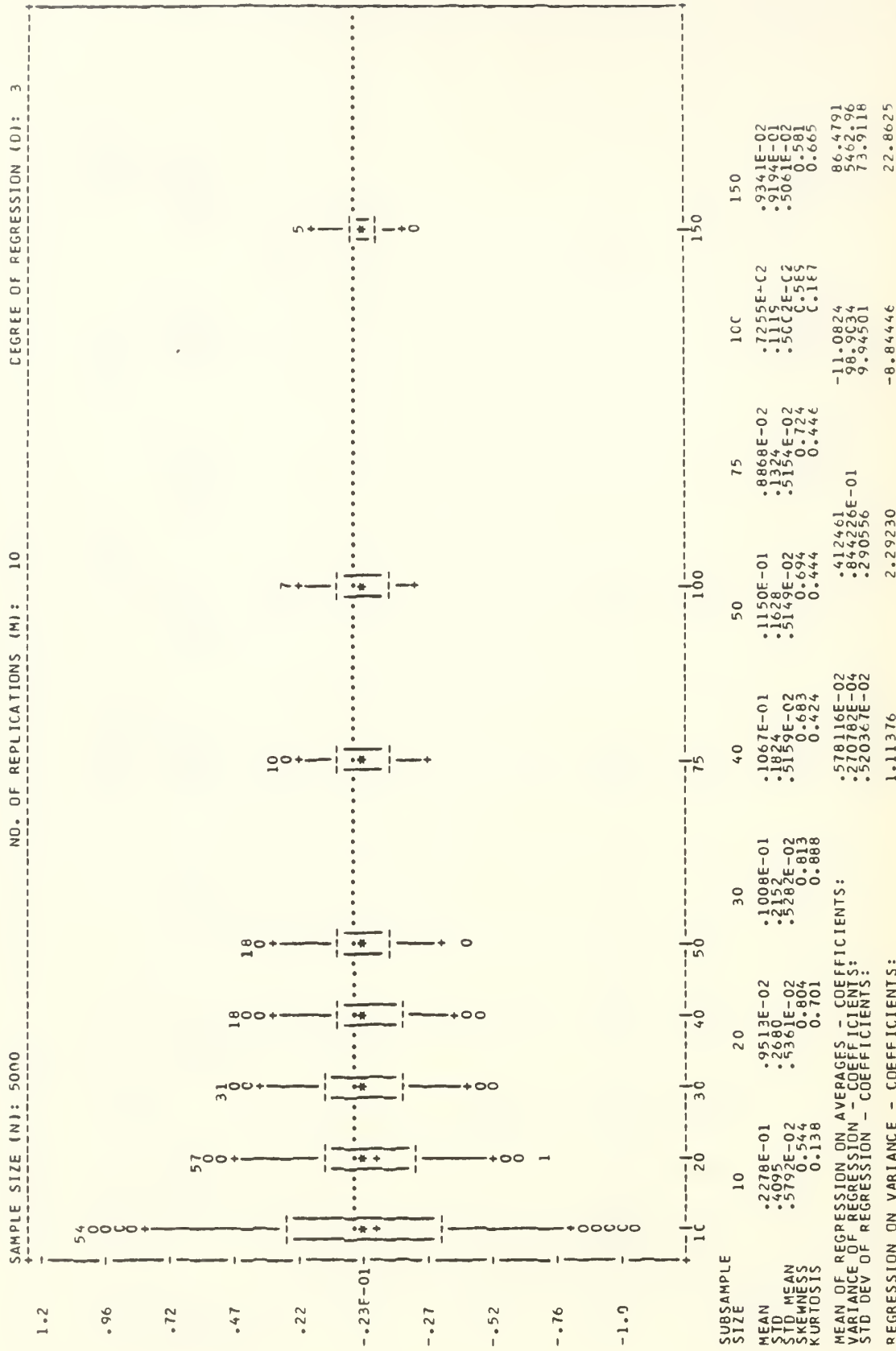
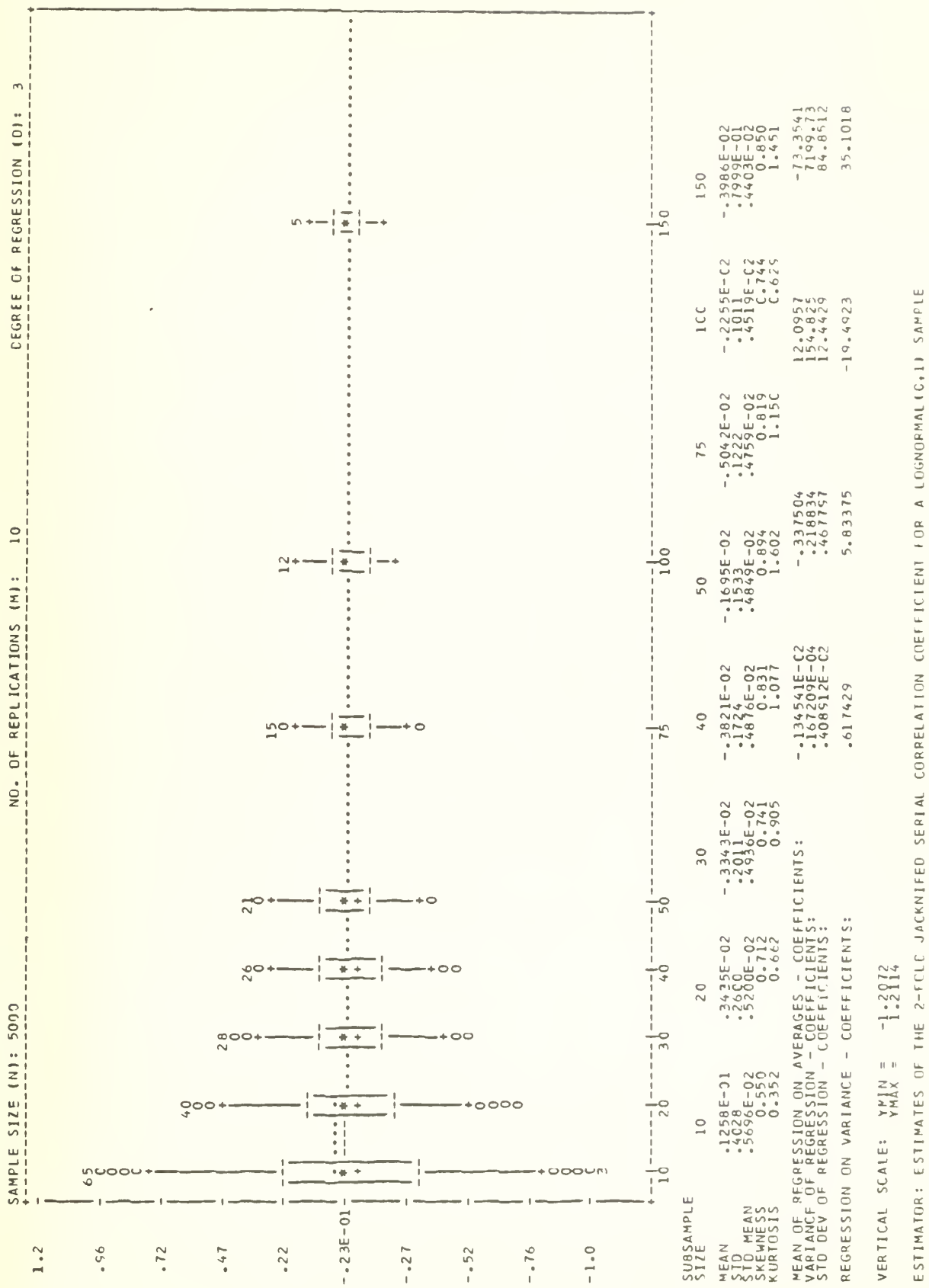


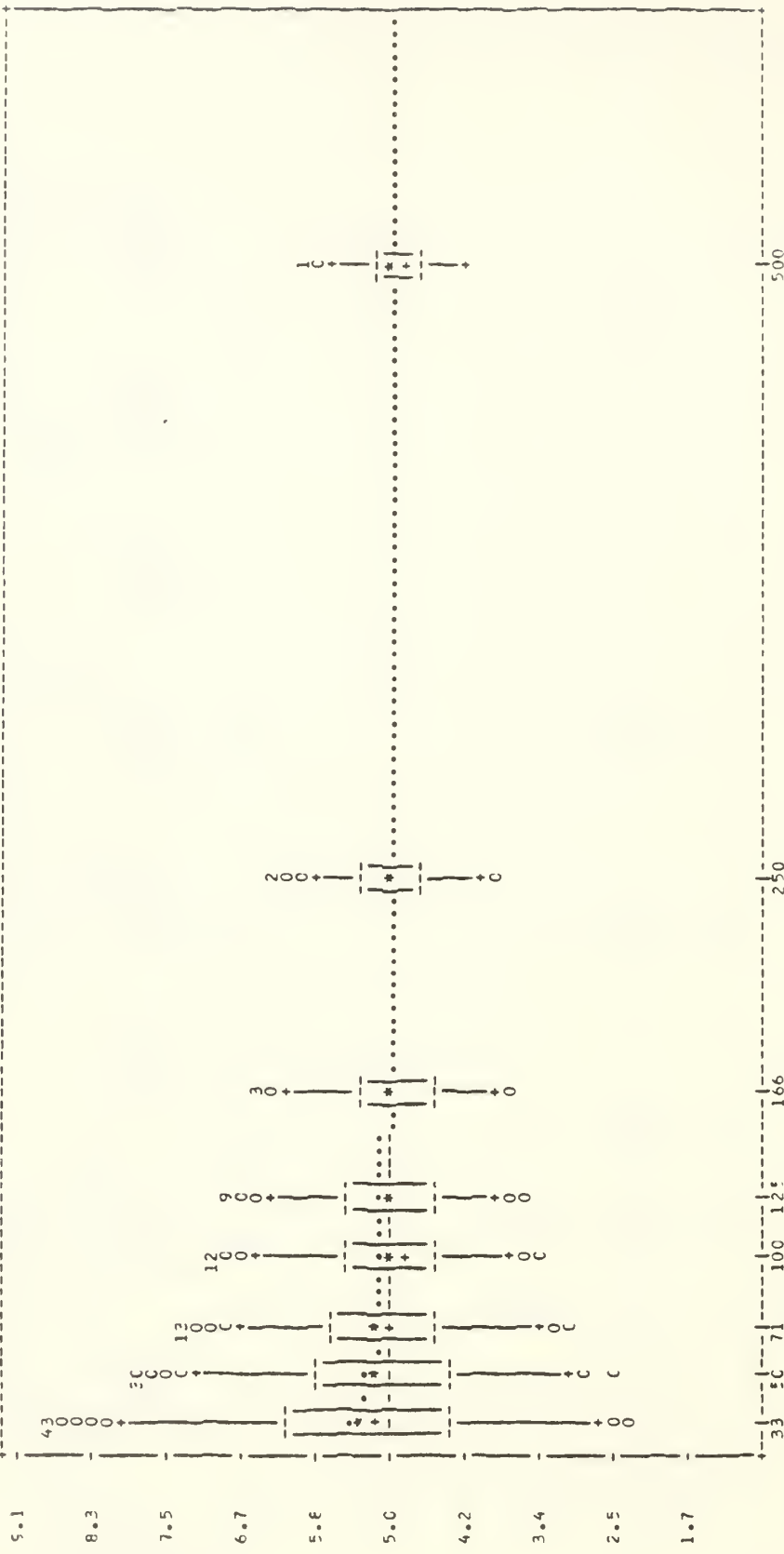
Figure 3(b)



SAMPLE SIZE (N): 2500

NO. OF REPLICATIONS (M): 20

DEGREE OF REGRESSION (D): 3



SAMPLE SIZE	33	50	71	100	125	166	250	500
MEAN	5.477	5.305	5.302	5.145	5.118	5.085	5.045	5.026
STD. MEAN	1.448	1.173	1.032	1.071	1.060	1.022	1.007	1.004
STD. DEVIATION	1.738E-01	1.355E-01	1.258E-01	1.343E-01	1.322E-01	1.324E-01	1.339E-01	1.365E-01
KURTOSIS	1.217	1.094	0.932	0.912	0.680	0.525	0.354	0.119
MEAN OF REGRESSION - COEFFICIENTS:	4.90635	12.5123	12.5123	12.5123	12.5123	12.5123	140.737	-2063.17
VARIANCE OF REGRESSION - COEFFICIENTS:	1.10504E-02	27.8280	27.8280	27.8280	27.8280	27.8280	180.94	5799.50E+08
STD. DEV. OF REGRESSION - COEFFICIENTS:	3.32424E-01	6.7684	6.7684	6.7684	6.7684	6.7684	351.610	7615.39
REGRESSION ON VARIANCE - COEFFICIENTS:	116.458	-1.681.96	-1.681.96	-1.681.96	-1.681.96	-1.681.96	14121.1	-34565.3

VERTICAL SCALE: YMIN = 1.0544
YMAX = 9.1332

ESTIMATOR: MAXIMUM LIKELIHOOD ESTIMATE OF THE SHAPE PARAMETER OF THE GAMMA DISTRIBUTION K=5.

Figure 4(a)

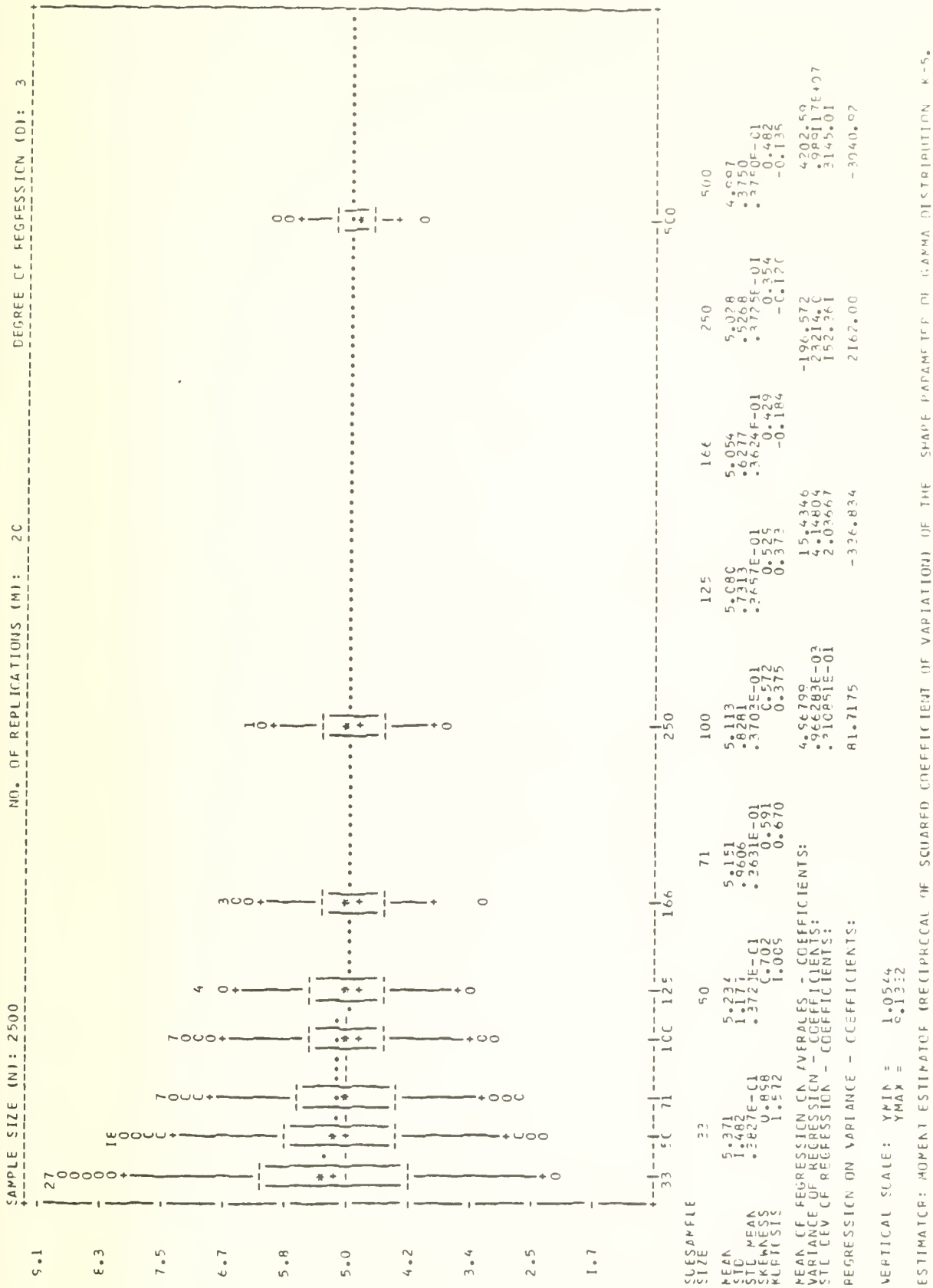


Figure 4(b)

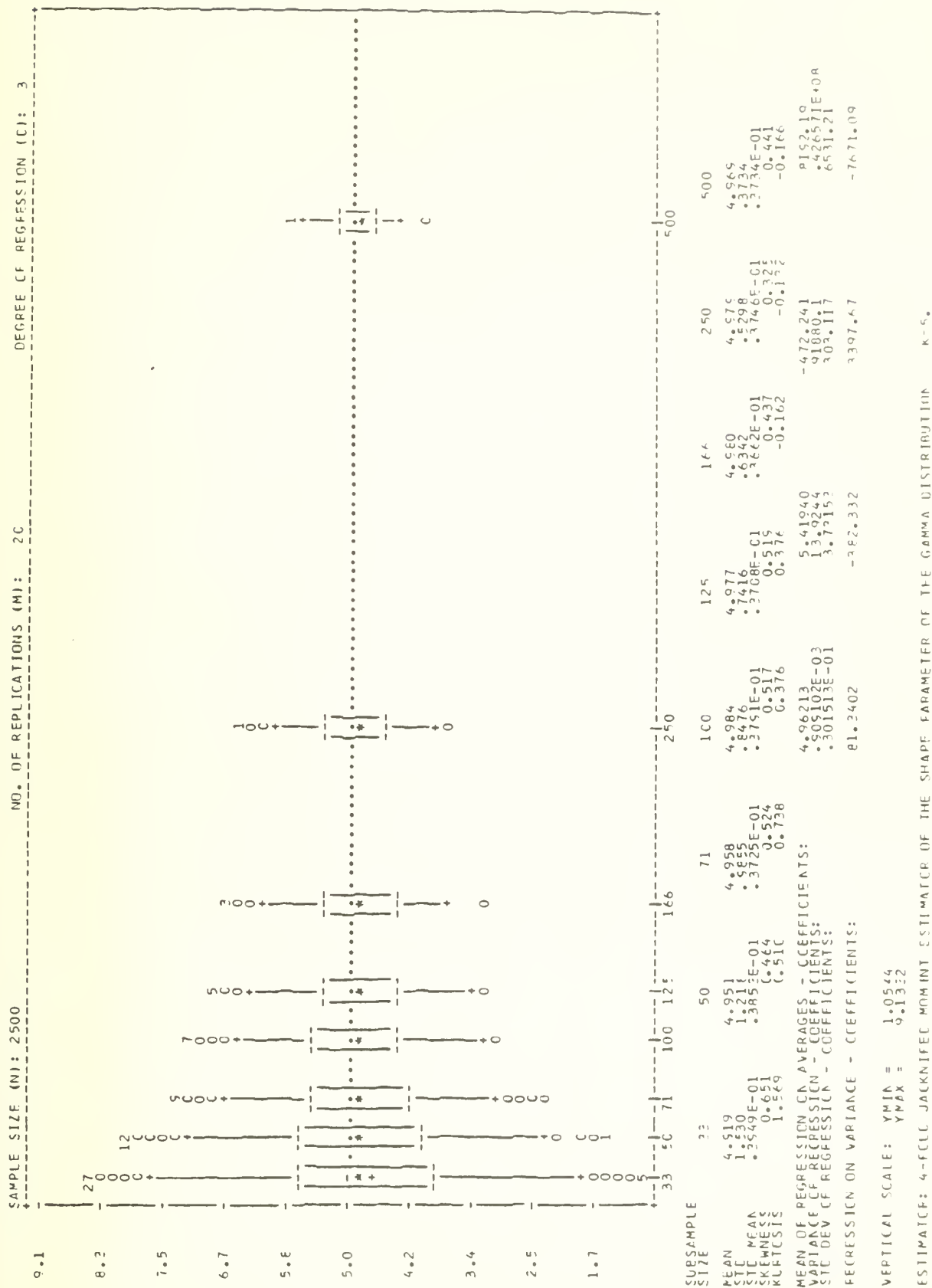
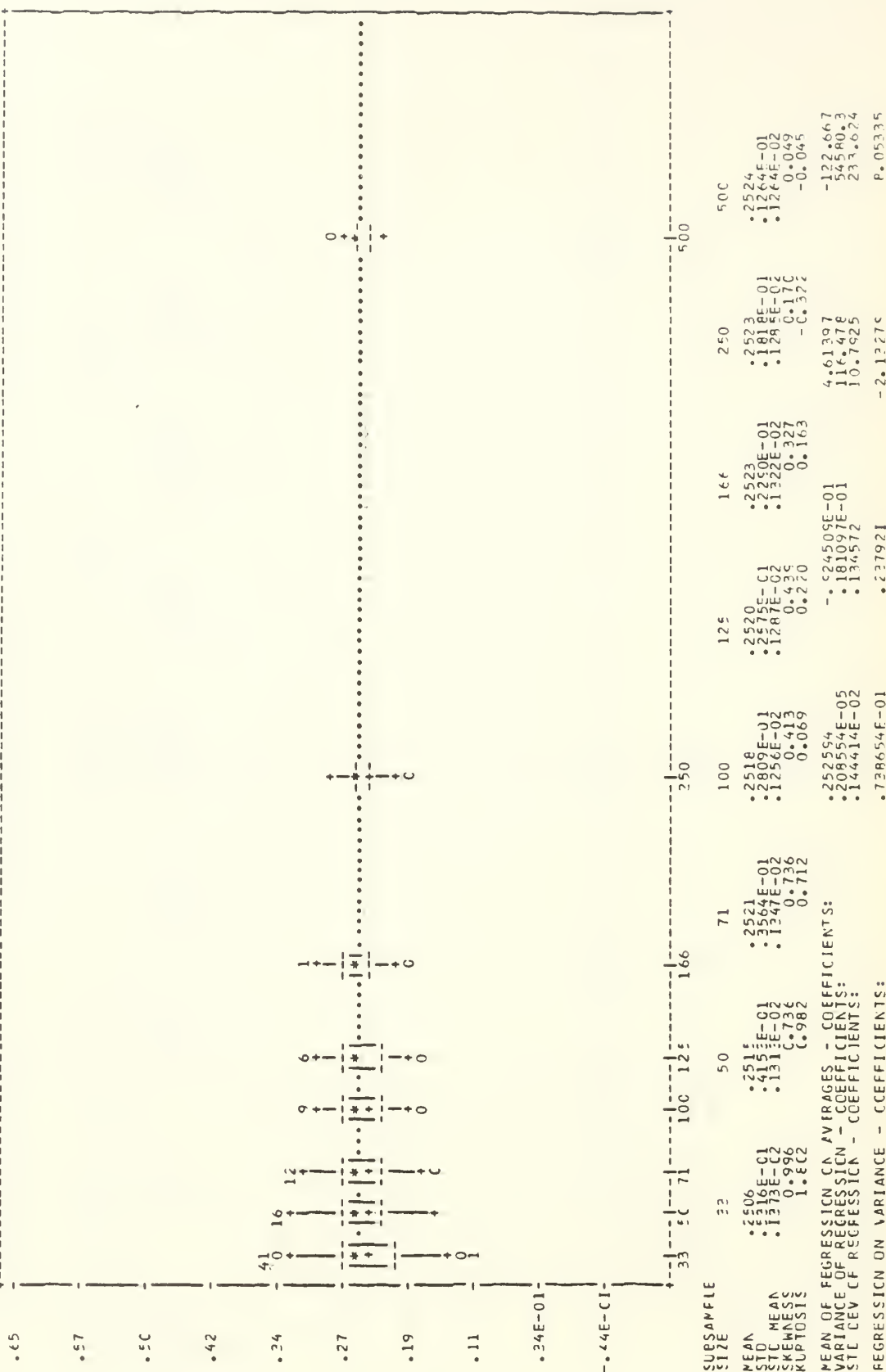


Figure 4(d)



VERTICAL SCALE: YMIN = -0.1055
YMAX = 0.6520

ESTIMATOR: 4-FOLD JACKKNIFE MAXIMUM LIKELIHOOD ESTIMATE OF THE SHAPE PARAMETER OF THE GAMMA DISTRIBUTION K=25

Figure 5(c)

DISTRIBUTION LIST

	NO. OF COPIES
Library, Code 0142 Naval Postgraduate School Monterey, CA 93940	4
Dean of Research Code 012A Naval Postgraduate School Monterey, CA 93940	1
Library, Code 55 Naval Postgraduate School Monterey, CA 93940	2
Professor P. A. W. Lewis Code 55Lw Naval Postgraduate School Monterey, CA 93940	170

DUDLEY KNOX LIBRARY - RESEARCH REPORTS



5 6853 01060362 4

~~U208577~~